

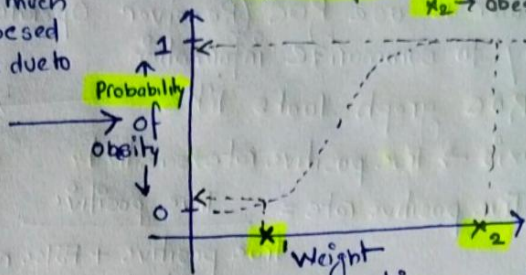
AUG and ROC

- Example - Whether the mouse is weight obese or not obese

This mouse doesn't weight that much but considered obese for its size (may be due to height)

This mouse have high weight still not considered obese because they might have lot of muscle

2 new sample - $x_1 \rightarrow$ Not obese
 $x_2 \rightarrow$ obese

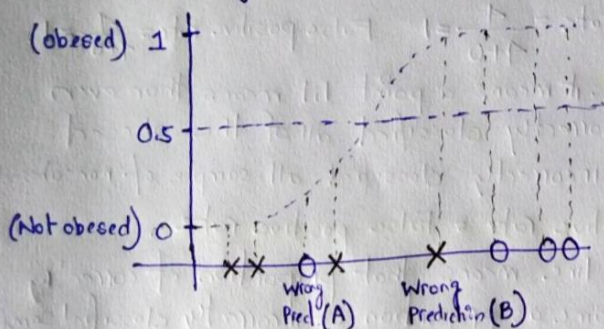


When we're doing logistic regression the y axis is converted to the probability of mouse is obese.

- However if we want to classify the mice as obese or not obese, then we need a way to turn the probabilities into classification.

- One way to classify is to set a threshold at 0.5. probability > 0.5 then obese.
probability < 0.5 then not obese.

- Suppose we get a four new set of mice who are NOT OBESE (X) and OBESE (O)

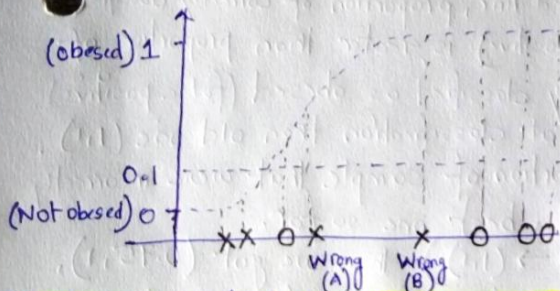


		Actual respectively	
		Obesed	Not obesed
Predicted	obesed	3	1
	not obesed	1	3

		Actual	
		Obesed	Not obesed
Predict	Obesed	TP	FP
	Not obesed	FN	TN

- suppose if is super important to correctly classify every obese sample, we can set the threshold to 0.1.

All 4 obesed (O) sample are now classified as obese.



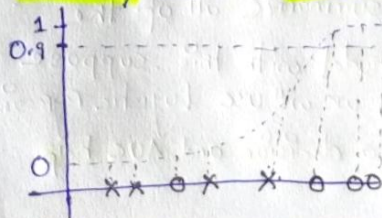
		Actual	
		obesed	Not obesed
Predicted	obesed	4	2
	Not obesed	0	2

→ Lower threshold would also reduced the number of False - Negative (FN) because all the obesed mice are correctly classified, and it would reduce True Negative (TN) because two of the mice that were not obese were incorrectly classified as obese.

→ So why we should decrease the threshold?

Suppose we are testing cancer or not cancer, then we should correctly classify who have cancer even if it means to decrease threshold.

- Otherway we can increase the threshold to 0.9 to correctly classify every non-obesed mice



		Actual	
		obesed	Not obesed
Predicted	Obesed	3	0
	Not obesed	1	4

→ We don't have any false positive

→ With higher threshold, does a better job classify sample as obese and not obese.

- Threshold can be set anything from 0 to 1.
- How do we determine which threshold is best?

Answer is use ROC (Receiver Operator Characteristics) graphs provide a simple way to summarize information.

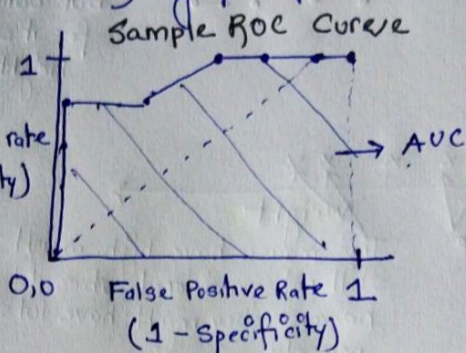
ROC graph looks like

Y axis \rightarrow True positive rate \rightarrow sensitivity

$$\text{True positive rate} = \frac{\text{True positive}}{\text{True positive} + \text{False negative}}$$

		Actual	
		obese	Not obese
Predicted	Is obese	TP	FP
	Not obese	FN	TN

True positive rate (Sensitivity)



True positive rate tells what proportion of obese samples are correctly classified.

X axis \rightarrow False positive rate $= 1 - \text{specificity} = \frac{\text{False positive}}{\text{False positive} + \text{True Negatives}}$

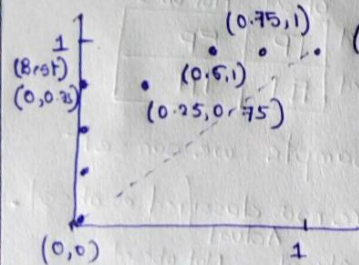
False positive rate tells proportion of non obese samples that were incorrectly classified and are false positive.

Suppose we set a threshold of 0.

		Actual	
		obese	Not obese
Predict	Obese	4	4
	Not obese	0	0

$$\text{True positive rate} = \frac{4}{4+0} = 1 \quad \text{False positive rate} = \frac{4}{4+0} = 1$$

Point be (1,1). It means a point 1,1 means that even though we correctly classified all of the obese samples, we incorrectly classified all samples of non-obese.



Line shows True positive rate = False positive rate.

Any point on this line means that proportion of correctly classified obese is same as proportion of incorrectly classified sample that not obese.

\rightarrow Suppose a new point is calculated (0.75,1)

(0.75,1) is to left of dotted line, proportion of correctly classified samples were obese (true positive) is greater than proportion of samples that are incorrectly classified as obese (false positive).

In short, new threshold which lead to (0.75,1) is better at classification than old one (1,1).

\rightarrow New point (0.5,1), this further decreased the proportion of sample that were incorrectly classified as obese (false positives). In other word, best one so far.

\rightarrow As we are increasing the threshold (started with 0 \rightarrow (1,1)) then we got (0.75,1), then (0.5,1) then (0.25,0.75), and (0,0.75). At (0,0.75) classify 75% of obese sample & 100% of the sample are not obese. In other word, this threshold result in no false positive.

\rightarrow The threshold is increased to 100% lead to (0,0) represent result is zero True positive & zero false positive.

\rightarrow If we connect the dot, ROC graph is generated. ROC graph summarize all of the confusion matrices that each threshold produced.

\rightarrow Area Under Curve (AUC) is 0.9. so we can compare two model with this. Suppose if AUC (logistic regression) is greater than AUC (random forest) then we will use logistic regression.

\rightarrow ROC curve make it easy to identify best threshold for making a decision and AUC help which method is better.